

Homework III

Due date: Sunday, Mar 17th 2024 at 11:59pm

Note: For written questions, you can turn in either a scanned copy of your handwritten answers or a PDF file of your answers. For programming questions, you need to submit your code and **also answer the asked questions in pdf**. Name the submission package as hw3_Lastname_Firstname.zip.

Problem 1: kNN Classifier (20 points)

You are asked to build a k -Nearest Neighbor (kNN) classifier. Here we use the heart data set. More information about the data can be found here¹. The data set is included in hw3_datasets.zip on Canvas. In the heart data folder, there are three files: “trainSet.txt”, “trainLabels.txt”, and “test.txt”. Each row of “trainSet.txt” corresponds to a data point whose class label is provided in the same row of “trainLabels.txt”. Each row of “testSet.txt” corresponds to a data point whose class label needs to be predicted. Code to load data is given in “p1.py”. You will train a classification model using “trainSet.txt” and “trainLabels.txt”, and use it to predict the class labels for the data points in “testSet.txt”.

1) Use the leave one out cross validation on the training data to select the best k among $\{1, 2, \dots, 10\}$. Report the averaged leave-one-out error (averaged over all training data points) for each $k \in \{1, 2, \dots, 10\}$.

2) Based on 1), use the best k to predict the class labels for test instances. You should also report the predicted labels for the testSet.

Problem 2: PCA (30 points)

You are asked to build a k -Nearest Neighbor (kNN) classifier based on dimensionality reduced data by PCA. The data set for evaluation is the gisette data set. More information about the data can be found here². The data set is included in hw3_datasets.zip on Canvas. The data is in the same format as that in Problem 1. You will train a classification model using “trainSet.txt” and “trainLabels.txt”, and use it to predict the class labels for the data points in “testSet.txt”. A demo code is given in “p2.py”, which includes data loading, model training, and visualization.

1) Train a k NN based on the original features. You should conduct cross-validation (of your choice) to select the best k . Describe the cross-validation approach, e.g., how many folds, how did you split the data. If you use a library for cross-validation, describe how the library does it. Report the best value of k under your cross-validation approach, and then report the testing accuracy corresponding to the best k .

¹[https://archive.ics.uci.edu/ml/datasets/Statlog+\(Heart\)](https://archive.ics.uci.edu/ml/datasets/Statlog+(Heart))

²<https://archive.ics.uci.edu/ml/datasets/Gisette>

2) Conduct PCA on the data and then learn a k NN model using the dimensionality reduced data. You should conduct cross-validation (of your choice) to select the best k and best d , where d is number of dimensions after PCA. Describe the cross-validation approach and report the best value of k and best value of d . Report the testing accuracy corresponding to the best k and best d .

3) Visualize the data distribution in the PCA projected space with $d=1$, $d=2$, and $d=3$, respectively. Include the visualization as figures in your pdf. Are there any visible clusters or patterns in the PCA plot? How does the PCA visualization change with different numbers of components?

Problem 3: Naïve Bayes for Text Classification (50 points)

Method In this problem, you are asked to implement a Naïve Bayes Classifier for text topic classification. Suppose there are m distinct words ($\{w_1, w_2, \dots, w_m\}$) in total within the training data. We represent a document \mathbf{x} by $\mathbf{x} = (x_1, x_2, \dots, x_m)$, where x_j is the occurrence frequency of word w_j in this document \mathbf{x} .

Given a collection of training documents $\mathbf{x}^1, \dots, \mathbf{x}^{n_k}$ in the class C_k , where $\mathbf{x}^i = (x_1^i, \dots, x_m^i)$, the probabilities of $p(w_j|C_k)$ can be estimated by MLE, i.e.,

$$p(w_j|C_k) = \frac{\sum_{i=1}^{n_k} x_j^i}{\sum_{j'=1}^m \sum_{i=1}^{n_k} x_{j'}^i}, \forall j, k.$$

To avoid the issue that some words may not appear in training documents of a certain class, the estimated probabilities are usually smoothed. One smoothing method is Laplace smoothing which computes $p(w_j|C_k)$ by

$$p(w_j|C_k) = \frac{\sum_{i=1}^{n_k} x_j^i + 1}{\sum_{j'=1}^m \sum_{i=1}^{n_k} x_{j'}^i + m}, \forall j, k.$$

Then, for a class C_k and a data point \mathbf{x} , we model $\Pr(\mathbf{x}|C_k)$ as

$$\Pr(\mathbf{x}|C_k) \propto \prod_{j=1}^m [p(w_j|C_k)]^{x_j}$$

where $p(w_j|C_k)$ stands for the probability of observing the word w_j in any document from the class C_k .

The log-likelihood for a data point $(\mathbf{x}, y = k)$ is given by

$$\log \Pr(\mathbf{x}, y = k) = \log \frac{\Pr(\mathbf{x}|C_k) \Pr(C_k)}{\Pr(\mathbf{x})} = \underbrace{\sum_{j=1}^m x_j \log p(w_j|C_k) + \log \Pr(C_k)}_{f_k(\mathbf{x})} + \text{const}$$

where $\Pr(C_k)$ can be estimated by $\Pr(C_k) = \frac{n_k}{\sum_{k=1}^K n_k}$ and $f_k(\mathbf{x})$ can be considered as a prediction score of \mathbf{x} for the k -th class. The class label of a test document \mathbf{x} can be predicted by $k^* = \arg \max_{1 \leq k \leq K} f_k(\mathbf{x})$.

Dataset The data set for training and evaluation is the 20NewsGroup data, which is included in the provided zip file. You will find six **text** files in this data set: train.data, train.label, train.map, test.data, test.label, and test.map, where the first three files are for training data and the last three files are for testing data. In the train.data file, you will find the word histograms of all documents; each row is a tuple of format (document-id, word-id, word-occurrence-frequency). The class labels of training documents can be found in train.label with the order corresponding to training documents' id, and the topic of each class can be found in train.map. Similarly, the word histograms and the class assignments of test documents can be found in test.data and test.label, respectively. Code to load data is given in “p1.py”.

Requirements In this problem, you need to:

- 1) Build a Naïve Bayes classifier with the Laplace smoothing using the training data.
- 2) Apply the learned classifier to predict the class labels for the test documents. Report the classification accuracy over the test documents (i.e., the proportion of test documents that are classified correctly). And also submit a file that contains the predicted labels of test documents corresponding to the original order.